

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**Methods and Systems for Adaptive Delivery of
Multimedia Contents**

Inventor(s):
Yudong Yang
Hong-Jiang Zhang

ATTORNEY'S DOCKET NO. MS1-871us

TO: 6545550

TECHNICAL FIELD

This invention relates to methods and systems for adaptive delivery of multimedia contents.

BACKGROUND

Today we live in a diversified world with a seemingly infinite number of diverse and different sources of on-line information. Typically, this information is accessed via a network such as the Internet. As the Internet and, more generally, computing evolve, people are beginning to become accustomed to and demanding better access to different types of electronically available information.

Against the backdrop of the diverse and different sources of electronic information is the wide range of devices that are connected being used to access the information. For example, people can now typically access network-accessible information using personal computers, handheld computers, personal digital assistants and the like. The situations encountered by individuals attempting to access this diverse collection of information using an ever growing collection of computing devices differs from user session to user session.

For example, some contents may be authored and suitable for use with a certain type of browser. Yet, the contents may not be suitably acceptable for use in other situations. For example, web contents that are authored for use in connection with a browser installed on a personal computer may not be suitable for display on a small handheld device for reasons not the least of which include the size disparity between the different devices' displays.

One possible solution for this problem is to author the same contents so that they reside in different forms that are suitable for all of the different situations that

1 might be encountered. While this is theoretically possible, the solution is
2 practically infeasible due to the time and expense involved.

3 One area of promise is in the area of so-called adaptive content delivery.
4 One goal of adaptive content delivery is to have content that is readily or easily
5 adaptable to different computing environments.

6 Early commercial applications focused on providing faster web page
7 downloads for narrow bandwidth connected users (such as dialup and mobile
8 access). Most of the applications accelerated downloads by simply reducing the
9 sizes of embedded image files using aggressive lossy compression schemes. The
10 cost of this solution is lower quality, which is highly undesirable from a customer
11 service standpoint. Some schemes also supported lossless text compression to
12 reduce the transmission time of web pages.

13 Some companies such as ProxyNet (based on TranSend technology),
14 SpyGlass, and OnlineAnywhere provide proxies or servers that can adjust web
15 pages to fit the display of smaller devices. Their technologies, however, are based
16 on heuristic rules and customized content filters that are designed for specific
17 websites and are used to extract the most important contents from these web
18 pages. Thus, these solutions tend to be rigid and inflexible.

19 Accordingly, this invention arose out of concerns associated with providing
20 adaptive systems and methods for efficient and flexible content delivery.

21 22 **SUMMARY**

23 Methods and systems for generic adaptive multimedia content delivery are
24 described. In one embodiment, a novel framework features an abstract content
25 model and an abstract adaptive delivery decision engine. The abstract content

1 model recognizes important aspects of contents while hiding their physical details
2 from other parts of the framework. The decision engine then makes content
3 adaptation plans based on the abstracted model of the contents and needs little
4 knowledge of any physical details of the actual contents. Thus, under the same
5 framework, adaptive delivery of generic contents is possible.

6 7 **BRIEF DESCRIPTION OF THE DRAWINGS**

8 Fig. 1 is an exemplary computer system that can be utilized in accordance
9 with one or more embodiments.

10 Fig. 2 is a block diagram of an exemplary adaptive content delivery system
11 in accordance with one embodiment.

12 Fig. 2a is a flow diagram that describes steps in a method in accordance
13 with one embodiment.

14 Fig. 2b is a flow diagram that describes steps in a method in accordance
15 with another embodiment.

16 Fig. 3 is a block diagram that illustrates an exemplary abstract content
17 representation structure in accordance with one embodiment.

18 Fig. 4 is a diagram that illustrates transfers of node statuses.

19 Fig. 5 is a diagram that illustrates an exemplary abstract content
20 representation structure of typical news content combined with possible transcoder
21 capabilities.

22 Fig. 6 is a diagram that illustrates an exemplary abstract content
23 representation structure of an MPEG I video sequence "IPBBIBBPBB..."

24 Fig. 7 is a graph of the unified QoS factor of the Fig. 6 MPEG I abstract
25 content representation structure nodes.

1 Fig. 8 contains graphs that illustrate aspects of adaptive delivery of
2 VBR/CBR MPEG I bitstreams.

3 Fig. 9 contains graphs that illustrate aspects of bandwidth versus frame
4 types, in accordance with an example that is given in the text.

5 6 **DETAILED DESCRIPTION**

7 **Overview**

8 Adaptive content delivery systems and methods are described. Efficiency
9 and flexibility are promoted through a novel solution to generic adaptive
10 multimedia content delivery. Described embodiments are based on an abstract
11 content model that captures important or critical structures and attributes of
12 contents. Contents are modeled as hierarchical directional graphs. Nodes on
13 graphs represent elements of contents. The concept of an "edge" is introduced.
14 Edges define logical relationships between these elements. By finding optimized
15 sub-graphs on these graphs under some constraints, optimized plans for adaptive
16 content delivery can be made. With the help of the abstract content model,
17 optimization procedures for many different types of contents can be standardized.
18 Accordingly different types of contents can be treated equally under this
19 framework.

20 21 **Exemplary Computer Environment**

22 The various components and functionality described herein can be
23 implemented by various computers. Fig. 1 shows components of a typical
24 example of such a computer, referred to by reference numeral 100. The
25 components shown in Fig. 1 are only examples, and are not intended to suggest

any limitation as to the scope of the claimed subject matter; the claimed subject matter is not necessarily dependent on the features shown in Fig. 1.

Generally, various different general purpose or special purpose computing system configurations can be used. Examples of well known computing systems, environments, and/or configurations that may be suitable for use with the invention include, but are not limited to, personal computers, server computers, hand-held or laptop devices, multiprocessor systems, microprocessor-based systems, set top boxes, programmable consumer electronics, network PCs, minicomputers, mainframe computers, distributed computing environments that include any of the above systems or devices, and the like.

The functionality of the computers is embodied in many cases by computer-executable instructions, such as program modules, that are executed by the computers. Generally, program modules include routines, programs, objects, components, data structures, etc. that perform particular tasks or implement particular abstract data types. Tasks might also be performed by remote processing devices that are linked through a communications network. In a distributed computing environment, program modules may be located in both local and remote computer storage media.

The instructions and/or program modules are stored at different times in the various computer-readable media that are either part of the computer or that can be read by the computer. Programs are typically distributed, for example, on floppy disks, CD-ROMs, DVD, or some form of communication media such as a modulated signal. From there, they are installed or loaded into the secondary memory of a computer. At execution, they are loaded at least partially into the computer's primary electronic memory. The invention described herein includes

these and other various types of computer-readable media when such media contain instructions programs, and/or modules for implementing the steps described below in conjunction with a microprocessor or other data processors. The invention also includes the computer itself when programmed according to the methods and techniques described below.

For purposes of illustration, programs and other executable program components such as the operating system are illustrated herein as discrete blocks, although it is recognized that such programs and components reside at various times in different storage components of the computer, and are executed by the data processor(s) of the computer.

With reference to Fig. 1, the components of computer 100 may include, but are not limited to, a processing unit 120, a system memory 130, and a system bus 121 that couples various system components including the system memory to the processing unit 120. The system bus 121 may be any of several types of bus structures including a memory bus or memory controller, a peripheral bus, and a local bus using any of a variety of bus architectures. By way of example, and not limitation, such architectures include Industry Standard Architecture (ISA) bus, Micro Channel Architecture (MCA) bus, Enhanced ISA (EISAA) bus, Video Electronics Standards Association (VESA) local bus, and Peripheral Component Interconnect (PCI) bus also known as the Mezzanine bus.

Computer 100 typically includes a variety of computer-readable media. Computer-readable media can be any available media that can be accessed by computer 100 and includes both volatile and nonvolatile media, removable and non-removable media. By way of example, and not limitation, computer-readable media may comprise computer storage media and communication media.

“Computer storage media” includes both volatile and nonvolatile, removable and non-removable media implemented in any method or technology for storage of information such as computer-readable instructions, data structures, program modules, or other data. Computer storage media includes, but is not limited to, RAM, ROM, EEPROM, flash memory or other memory technology, CD-ROM, digital versatile disks (DVD) or other optical disk storage, magnetic cassettes, magnetic tape, magnetic disk storage or other magnetic storage devices, or any other medium which can be used to store the desired information and which can be accessed by computer 100. Communication media typically embodies computer-readable instructions, data structures, program modules or other data in a modulated data signal such as a carrier wave or other transport mechanism and includes any information delivery media. The term “modulated data signal” means a signal that has one or more of its characteristics set or changed in such a manner as to encode information in the signal. By way of example, and not limitation, communication media includes wired media such as a wired network or direct-wired connection and wireless media such as acoustic, RF, infrared and other wireless media. Combinations of any of the above should also be included within the scope of computer readable media.

The system memory 130 includes computer storage media in the form of volatile and/or nonvolatile memory such as read only memory (ROM) 131 and random access memory (RAM) 132. A basic input/output system 133 (BIOS), containing the basic routines that help to transfer information between elements within computer 100, such as during start-up, is typically stored in ROM 131. RAM 132 typically contains data and/or program modules that are immediately accessible to and/or presently being operated on by processing unit 120. By way

of example, and not limitation, Fig. 1 illustrates operating system 134, application programs 135, other program modules 136, and program data 137.

The computer 100 may also include other removable/non-removable, volatile/nonvolatile computer storage media. By way of example only, Fig. 1 illustrates a hard disk drive 141 that reads from or writes to non-removable, nonvolatile magnetic media, a magnetic disk drive 151 that reads from or writes to a removable, nonvolatile magnetic disk 152, and an optical disk drive 155 that reads from or writes to a removable, nonvolatile optical disk 156 such as a CD ROM or other optical media. Other removable/non-removable, volatile/nonvolatile computer storage media that can be used in the exemplary operating environment include, but are not limited to, magnetic tape cassettes, flash memory cards, digital versatile disks, digital video tape, solid state RAM, solid state ROM, and the like. The hard disk drive 141 is typically connected to the system bus 121 through an non-removable memory interface such as interface 140, and magnetic disk drive 151 and optical disk drive 155 are typically connected to the system bus 121 by a removable memory interface such as interface 150.

The drives and their associated computer storage media discussed above and illustrated in Fig. 1 provide storage of computer-readable instructions, data structures, program modules, and other data for computer 100. In Fig. 1, for example, hard disk drive 141 is illustrated as storing operating system 144, application programs 145, other program modules 146, and program data 147. Note that these components can either be the same as or different from operating system 134, application programs 135, other program modules 136, and program data 137. Operating system 144, application programs 145, other program

1 modules 146, and program data 147 are given different numbers here to illustrate
2 that, at a minimum, they are different copies. A user may enter commands and
3 information into the computer 100 through input devices such as a keyboard 162
4 and pointing device 161, commonly referred to as a mouse, trackball, or touch
5 pad. Other input devices (not shown) may include a microphone, joystick, game
6 pad, satellite dish, scanner, or the like. These and other input devices are often
7 connected to the processing unit 120 through a user input interface 160 that is
8 coupled to the system bus, but may be connected by other interface and bus
9 structures, such as a parallel port, game port, or a universal serial bus (USB). A
10 monitor 191 or other type of display device is also connected to the system bus
11 121 via an interface, such as a video interface 190. In addition to the monitor,
12 computers may also include other peripheral output devices such as speakers 197
13 and printer 196, which may be connected through an output peripheral interface
14 195.

15 The computer may operate in a networked environment using logical
16 connections to one or more remote computers, such as a remote computer 180.
17 The remote computer 180 may be a personal computer, a server, a router, a
18 network PC, a peer device or other common network node, and typically includes
19 many or all of the elements described above relative to computer 100. The logical
20 connections depicted in Fig. 1 include a local area network (LAN) 171 and a wide
21 area network (WAN) 173, but may also include other networks. Such networking
22 environments are commonplace in offices, enterprise-wide computer networks,
23 intranets, and the Internet.

24 When used in a LAN networking environment, the computer 100 is
25 connected to the LAN 171 through a network interface or adapter 170. When used

1 in a WAN networking environment, the computer 100 typically includes a modem
2 172 or other means for establishing communications over the WAN 173, such as
3 the Internet. The modem 172, which may be internal or external, may be
4 connected to the system bus 121 via the user input interface 160, or other
5 appropriate mechanism. In a networked environment, program modules depicted
6 relative to the computer 100, or portions thereof, may be stored in the remote
7 memory storage device. By way of example, and not limitation, Fig. 1 illustrates
8 remote application programs 185 as residing on computer 180. It will be
9 appreciated that the network connections shown are exemplary and other means of
10 establishing a communications link between the computers may be used.

11 12 **Exemplary Embodiment**

13 In the embodiment about to be described, the inventive adaptive content
14 delivery system works as an extended content processor of traditional multimedia
15 content servers. Accordingly, upon receipt of a content request, the server fetches
16 original multimedia contents from a content source and passes them to the
17 adaptive delivery system. Adaptation results are then sent as a response to the
18 request.

19 Fig. 2 shows an adaptation system 200 in accordance with one
20 embodiment. System 200 includes a content host 202 that acts as a front-end
21 between physical multimedia contents and the rest of the system. The content host
22 comprises a content parser 204 and a content mapper 206. Content host 202
23 processes the original multimedia contents to provide an abstract content model
24 which is provided to a decision engine 208 for processing. In this example, the
25 abstract content model can be the only data of which other system components

(e.g. the decision engine) are aware. Physical details of the original multimedia content are thus hidden by the abstract content model. Other input parameters, such as network characteristics and client capabilities, are modeled as resources (i.e. resource model 210) or preference factors (i.e. preference model 212). The decision engine 208 then makes best-fit delivery plans based on these parameters and the abstract content model. Each of these components is explored in more detail below.

Content Host

As the illustrated and described content host 202 is a front-end between physical contents and other system components, it can be used to manipulate contents based on the abstract content model. In addition, the content host 202 also defines a common set of application program interfaces or APIs for retrieving extended properties of the abstract content model. Exemplary APIs are given at the end of this document. Although the remaining components of the system are content independent, content host 202 itself is dependent on content types. Thus, different media types may desire different implementations of the content host. As a common basis, however, the content host 202 should comprise two sub-modules: content parser 204 and content mapper 206.

In the illustrated and described embodiment, content parser 204 scans input contents and constructs corresponding abstract content model representations either online or offline. Different formats of the same content and capabilities of supported transcoders can also be abstracted into the same model during this process. Specific examples of how this can be done are given below. Characteristics of one suitable technique for implementing the content parser are

described in U.S. Patent Application Serial No. 09/893,335, entitled "Function-based Object Model for Use in WebSite Adaptation" filed on June 26, 2001, the disclosure of which is incorporated by reference herein.

Content mapper 206 functions in a manner that is opposite of the way content parser 204 functions. That is, content mapper 206 converts abstract content model representations back to physical contents. Real-time-capable content transcoders may also be called at this stage to generate desired results.

Decision Engine

Decision engine 208 provides functionality for making content adaptation plans. In the illustrated example, decision engine 208 selects appropriate contents that achieve maximum total QoS (i.e. quality of service) values according to current resource constraints and preference factors (as provided by resource model 210 and preference model 212). Based on the abstract content model, this problem is solved by finding optimized sub-graphs of the abstract content model that maximize QoS values under resource constraints. Details of an exemplary content optimization procedure are covered in the section entitled "Content Optimization" below.

Resource and Preference Models

Input parameters, such as network characteristics and client capabilities, are modeled as resources or preference factors. Resources are used as constraints while the decision engine is looking for the best delivering plans. Preference factors are used to alter QoS factors of the abstract content models. For dynamically changing parameters, such as network characteristics, these models

should be able to predict future values since the decision engine may use forward-looking algorithms. More information is provided on this topic in the section entitled “A Simple Sub-optimization Algorithm” below.

Caching

Although a caching stage is not explicitly included Fig. 2, caches can play an important role in real adaptive content delivery systems. Complex transcoding processes can be avoided if the needed results are in cache. Partial plans of delivery can also be saved in the cache and reused under nearly identical content request situations. These savings will decrease the server-side resource requirements and thus better quality of service (QoS) can be achieved when server-side resources become the bottleneck of the delivery plans.

Fig. 2a is a flow diagram that describes steps in a method in accordance with one embodiment. The steps can be implemented in any suitable hardware, software, firmware or combination thereof. In the present example, the steps are implemented in software. In addition, any suitable software architecture can be utilized to implement the steps about to be described. One exemplary software architecture is shown and described in connection with Fig. 2 above. It should be appreciated, however, that Fig. 2 shows but one exemplary software architecture and should not be construed to limit application of the claimed subject matter except where so specifically recited. In the example about to be described, the processing steps reflect what can be considered as an “on-line” mode. By “on-line” is meant that the processing can take place responsive to a content request. Aspects of the described processing can, however, take place in an “off-line” mode where contents are pre-processed, for example, prior to receiving a specific

request for the contents. For example, when content is received by a server for storage, the content can be pre-processed to build the abstract content model. Subsequently, when any requests are received, the server can simply retrieve the abstract content model and select an appropriate delivery plan. Such is explored in more detail in Fig. 2b.

Step 250 receives a content request. This step can be implemented responsive to a client device sending such a request. Step 252 retrieves the requested content from a content source and step 254 parses the content and builds an abstract content model. An exemplary abstract content model is described below in more detail. Step 256 processes the abstract content model to select an optimal delivery plan. Examples of how this can be done are described below. Step 258 then processes the abstract content model to provide deliverable content in accordance with the delivery plan. Step 260 then delivers the content to the content requester.

Fig. 2b is a flow diagram that describes steps in a method in accordance with another embodiment. The steps can be implemented in any suitable hardware, software, firmware or combination thereof. In the present example, the steps are implemented in software. In addition, any suitable software architecture can be utilized to implement the steps about to be described. The processing about to be described can pertain to the "off-line" mode mentioned above.

Step 262 receives content. This step can be implemented in any suitable way. For example, this step can be implemented when a server receives content that it is to store for future content requests. Step 264 parses the content and builds an abstract content model. An exemplary abstract content model is described below in more detail. Step 266 processes the abstract content model to

contents, and edges between the nodes represent relationships between these components. In the discussion that follows, definitions of these data models and their basic attributes are given. Then, a discussion of the details of nodes and edges is presented.

For the discussion that follows, the reader is referred to Fig. 3 which provides an illustration of an exemplary abstract content representation structure 300 in accordance with one embodiment.

In the illustrated structure 300, a *node* is an abstract representation of content or content structure. A Node is represented using a circle or square in the drawings of the data structure. In Fig. 3, five exemplary nodes are indicated at 302, 304, 306, 308, and 310. An *edge* is an abstract representation of relationships between the nodes. In this embodiment, edges can be divided into three different types. A *dependency edge* defines a logical dependency between nodes. In Fig. 3, a thin dashed line is used to represent dependency edges. Dependency edges are seen to extend between nodes 304, 310 and 302, 306. A *route edge* defines an ordered or hierarchical dependency between nodes. Thick solid lines are used to represent this kind of edge. In Fig. 3, a route edge extends between nodes 306, 310. A *mixed edge* is a mixture of a dependency edge and a route edge and can be considered as two edges separately. Thick dashed lines are used in Fig. 3 to represent this type of edge. A mixed edge can be seen to extend between nodes 304, 306 and 306, 308.

Abstract content representation structure 300 comprises a directional graph $G=(N, E)$ that satisfies the following condition (layered constraint), where N and E stand for “node” and “edge” sets of G :

- The node set N can be divided into several subsets as N_i , $i=1..m$ where $N = \bigcup_{i=1}^m N_i$ and $N_i \cap N_j = \Phi$,

- The node set N can be divided into several subsets as N_i ,
 $i=1..m$ where $N = \bigcup_{i=1}^m N_i$ and $N_i \cap N_j = \Phi$,
 $\forall i, j = 1..m, i \neq j$.
- $\forall a = (n_s, n_e) \in E$, there exists a pair of $1 \leq i < j \leq m$ such that $n_s \in N_i$
and $n_e \in N_j$.
- $\forall a = (n_s, n_e) \in E, s \neq e$.

From this definition one can also see that nodes in N_l do not have incoming edges and nodes in N_m do not have out going edges. In addition to these definitions, nodes and edges have several basic attributes as listed below.

Node and Edge Attributes

Node status defines the dynamic statuses of nodes during content delivery. In the illustrated and described embodiment, the node statuses include the following:

- ***Inactive***: The node is not yet a deliverable object because of an unsatisfied active condition (defined below).
- ***Activable***: The active condition of the node is satisfied but the node is not yet included in a delivery plan.
- ***Activated***: The node has been chosen in a delivery plan.
- ***Delivered***: The node has been delivered successfully to a content receiver.

- 1
- **Skipped:** The node is not delivered and will not be included in the
2 current plan due to some reason.
- 3

4 Fig. 4 illustrates some possible transfers between these statuses.
5 Specifically, an inactive node can become activable and vice versa. In addition, an
6 activated node can become inactive. Activable nodes can become activated or
7 skipped, and so on.

8 An *ignition edge* is defined as a dependency edge from a node that is
9 activated, delivered or skipped. In Fig. 3, which statuses are to be considered are
10 specified deliberately and marked as '+', '*' and '-' respectively, or no markers
11 are drawn if all three statuses are fine, i.e. if no action is to be taken. For example,
12 in Fig. 3, the edge between nodes 302 and 306 is an ignition edge if node 302 is
13 activated or delivered.

14 An *active condition* defines how the node becomes activable. In the
15 illustrated example, there are three conditions:

16

- **Automatic:** The node is automatically activable. One example is
17 the top most nodes (such as those in set N_I). These nodes are
18 considered as the potential starting point of delivery. In the
19 drawings in this document, automatically activable nodes are
20 designated mark with a "#".
- **OR condition:** The node is activable if at least one of its input edges
21 is an ignition edge. There is no marker for these kinds of nodes in
22 the drawings.
23
24
25

FIG. 1

- **AND condition:** The node is activable only if all its input edges are ignition edges. These nodes are designated with a “&” in the drawings.

An **input condition** defines when an activable node can be activated. In the present embodiment, except for automatic activable nodes, the only condition is that there exists an incoming route edge from an activated, delivered or skipped node. An **output behavior** defines how nodes on ends of outgoing route edges can be branched. In the illustrated and described embodiment, a “branch” comprises the operation of changing an activable node to an activated node. In this example, there are three branch operations:

- **Complementary branch:** All activable nodes can be branched. In the drawings in this document, there are no markers for these nodes.
- **Exclusive branch:** Only one from all activable nodes can be branched. In the drawings in this document, these nodes are designated with a “%”.
- **Tight branch:** All activable nodes must be branched. A “\$” is used to designate these nodes.

Nodes and Edges

As indicated above, nodes and edges are the basic elements of the described abstract content representation structure. Nodes and edges represent content objects and relationships between these objects. In addition to the basic attributes

introduced and discussed above, nodes and edges can have other extended attributes.

Nodes

As an abstraction of multimedia contents, a node can represent a piece of raw content like a picture, a paragraph of text, a video frame, or a chunk of bits in some coded contents. It can also work as a connection point of content structures where it does not represent any actual contents. In addition to its basic attributes, nodes can have some additional application related attributes.

A *value* (QoS factor) is defined as an increment to content quality when this node has been delivered. For a node $n \in N$, its value is represented by $V(n)$ or $V(n, t)$ if it is related to time.

A *resource factor* defines an amount of resources needed to deliver the node. For some types of resources r , one can represent corresponding resource factors of node n by $R(n, r)$.

Edges

Edges can represent any kind of relationship between node objects. Some possible edge meanings are listed below:

- *Semantic inclusion*: A node at the end of an edge is included in some semantic scope of the node at the beginning of the edge. E.g.: news stories (each as a node) of a specific category (a structural node), chapters of a book, etc.

- Fig. 5 shows an example that uses an abstract content representation structure to represent contents associated with typical news data and some transcoding scenarios. Starting from the topmost nodes, the news stories are classified or categorized as world, domestic, etc. Then in each category, there may be standalone stories and collections of specific topics that contain related stories (such as the node designated “Story 1”). A story can include different contents and these contents can be transcoded into different types, such as multiple

Fig. 5 shows an example that uses an abstract content representation structure to represent contents associated with typical news data and some transcoding scenarios. Starting from the topmost nodes, the news stories are classified or categorized as world, domestic, etc. Then in each category, there may be standalone stories and collections of specific topics that contain related stories (such as the node designated “Story 1”). A story can include different contents and these contents can be transcoded into different types, such as multiple

Fig. 5 shows an example that uses an abstract content representation structure to represent contents associated with typical news data and some transcoding scenarios. Starting from the topmost nodes, the news stories are classified or categorized as world, domestic, etc. Then in each category, there may be standalone stories and collections of specific topics that contain related stories (such as the node designated “Story 1”). A story can include different contents and these contents can be transcoded into different types, such as multiple

Content Optimization

Using the abstract content representation structure, optimized content can be considered as an optimal sub-graph of the corresponding structure. One target of optimization is to maximize preference-altered total QoS factors of abstract content representation structure nodes covered by the sub-graph under the constraints of the resources. In the section immediately below, some suggestions on choosing proper QoS factors are presented. Then a discussion of a simple bounded search algorithm as a near optimal solution to this content optimization problem will be presented. The algorithm is also used in a verification prototype that is introduced in the section entitled “Adaptive MPEG I Video Streaming” below.

Choosing Proper QoS Factors

The QoS factors of abstract content representation structure nodes play an important role in content optimization. This is because the decision engine 208 (Fig. 2) is programmed to decide which content is more suitable based on QoS factors. Accordingly, it can be advantageous for properly selected QoS value

1 definitions to at least conform to the following two principles. First, the QoS
2 value definitions should reflect the importance of the corresponding content.
3 Second, the QoS values of different nodes should be comparable.

4 In many cases, the first principle is easier to follow and conforming to the
5 second principle is usually not trivial. It may not be easy to tell which is more
6 important or meaningful as between two different contents. For example, there is
7 a famous saying that says "a picture is worth a thousand words". This might be
8 true in some cases, but not in others. In resource critical applications such as
9 mobile communication, text should be more preferable than images most of the
10 time. Thus, it is suggested that QoS definition choices be made on an application
11 specific basis.

12 A Simple Sub-Optimization Algorithm

13 In the illustrated and described embodiment, a bounded search algorithm is
14 adopted to find the near optimal solution of the content optimization problem. The
15 pseudo code listed below describes but one optimized adaptive content delivery
16 algorithm. The algorithm is a straightforward implementation of a deep-first
17 search.
18

```
19 (node, value) BestNode(ACRES, step, max_steps)
20 {
21     if step >= max_steps then return (NULL, 0);
22     BestCandi = NULL;
23     BestValue = 0;
24     Candi_Set = GetActivableNodes(ACRES);
25     Resource = GetFreeResource();
26     for each candi in Candi_Set do
27     {
28         resCandi = GetResource(candi);
29         if resCandi > Resource then next;
30         MarkActivated(candi);
31         ConsumeResource(resCandi);
32         valCandi = GetQoSValue(candi)
```



```

1         *GetPreference(candi)
2         + BestNode(ACRES, step+1, max_steps).value;
3         if valCandi > BestValue then
4         {
5             BestCandi = candi;
6             BestValue = valCandi;
7         }
8         FreeResource(resCandi);
9         MarkActivable(candi);
10        }
11        return (BestCandi, BestValue);
12    }
13
14    AdaptiveDelivery(ACRES, max_steps)
15    {
16        do
17        {
18            (node, value) = BestNode(ACRES, 0, max_steps);
19            if node != NULL then
20            {
21                MarkActivated(node);
22                ConsumeResource(GetResource(node));
23                Deliver(node);
24            }
25            UpdateSystemStatus();
26        }
27        while Not Meet_End_Condition();
28    }

```

The pseudo code starts from a candidate set of activable nodes and then tries to simulate following delivery plans by marking nodes on the path as activated temporarily. A back trace is then used to find other possible delivery plans. Finally, the best starting candidate is selected and delivered. Afterwards, system statuses, such as resources and preferences, are updated. The algorithm is then looped until the end condition is met. It will be appreciated that in some cases, dynamically changing factors, such as network resources, may have to be predicted during the search.

Because there are only limited search depths, the complexity of this algorithm is quite acceptable. However, the algorithm may become complex in some cases. Pruning of search branches is not currently done because node QoS

1 factors may depend on resource consumption of previously selected nodes and
2 thus may change dynamically. Under such a situation, historical records are not
3 reusable and searching cannot be accelerated by pruning. In order to benefit from
4 pruning, modifications can be made. For example, one way to benefit from
5 pruning is to replace continuous values with approximate discrete ones (as time,
6 QoS, resource).

7 Adaptive MPEG I Video Streaming Example

8 As a verification and example of the inventive framework, a simple
9 adaptive MPEG I video streaming application was implemented and based on the
10 content optimization algorithm and above-described abstract content
11 representation structure. One goal of this application is to allow smoothed
12 playback of MPEG I video even when transmission bandwidth is less than that
13 which the original video bitstream requires and/or the client-side buffer size is
14 limited.

15 In the discussion that follows, we start from an analysis of the situation and
16 then we will build an abstract content representation structure of the MPEG I
17 video bit stream that is used in our adaptive delivery verification prototype. After
18 that, a brief introduction will be given to our implementation's architecture.
19 Experimental results are discussed later.

20 Abstract Content Representation Structures and Adaptation Methods in 21 Detail

22 Since an adaptable abstract content representation of content is the basis of
23 the inventive approach, this discussion starts by analyzing adaptation schemes of
24 MPEG I video and then constructs the abstract content representation model that
25 enables the adaptation scheme based on the analysis.

1 As it was designed, MPEG I video bitstreams do not support scalable
2 delivery. Network bandwidth must be large enough to enable smooth playbacks in
3 normal cases. Transcoders are required if network bandwidth is less than the
4 original video requires. However, online transcoding of MPEG video bitstreams
5 is computing intensive and does not suit video streaming applications where
6 multiple contents/connections need to be supported simultaneously. For this
7 reason, a simpler adaptation approach is chosen by selectively replacing/skipping
8 encoded video frames. This approach is very efficient and the video bitrate is
9 reduced at the cost of lower frame rates instead of PSNR losses resulting from
10 normal transcoding.

11 There are three kinds of encoded video frames in MPEG I video bitstreams
12 -- Intra coded, forward Prediction coded and Bi-directional prediction coded.
13 These frame types provide a trade-off between compression efficiency and
14 playback requirements (as seek and error recovery). Several encoded frames form
15 a GOP (group of picture) and temporal references of frames are defined relatively
16 within GOPs. Beside a sequence header that defines essential attributes, video
17 bitstreams are typically a concatenation of GOPs. I frames can be decoded at
18 anytime, but decoding of P and B frames depends on decoded reference frames.
19 In other words, there are relationships that exist as inter-frame dependencies and
20 temporal orders. As a result, skipping I and P frames will cause P and B frames
21 that follow not to be decodable. Skipping B frames has no impact on other
22 frames.

23 Beside these inter-frame dependencies, frame timing is another issue that is
24 addressed during content adaptation. Although video bitstreams with some
25 skipped frames can be decoded without any problems, the frame timing is

changed. As a remedy, escape-coded frames are used instead of skipping where PB frames should be skipped. Thus frame timing is kept unchanged during playback. The escape-coded frames are MPEG I coded frames too and they stand for nothing changed to the reference frame (or one of the reference frames if the frame type is B). According to MPEG I syntax, all macroblocks of a frame must be covered by non-overlapped slices, and a slice must start and end by coded macroblocks. Thus, the minimal escape coded frame must consist of at least two empty coded macroblocks (top-left and bottom-right) and address skip codes for all other macroblocks between them. As a result, the minimal size of an escape coded CIF frame (P or B) is 32 bytes which is minor, if compared to that of normally coded frames which are at least several kilobytes.

In accordance with the above discussion, Fig. 6 shows an exemplary abstract content representation model of a typical MPEG I video bitstream. Each node in the figure represents data of a coded frame as "I", "P", "B". Additionally, the designations of "a" or "b" represent data of escape-coded frames for P or B frames respectively. The temporal reference of each frame is shown as a number before the frame type. Data of the sequence header and GOP headers are not shown in this figure due to limited space.

Each node in this model has attributes including data size, expected time to be decoded (TTD) relative to the starting time of delivery, and QoS factor. The QoS factor is defined dynamically according to the current time and TTD. In this example, this value is assigned based on the following heuristics:

- The value of frame data depends on frame types and effects on decoding of succeeding frames. $I > P > B > \text{escape-coded}$.

- The value of a frame should be maximized when it is available for the decoder just on or before its TTD and decrease to zero when it is not delivered after TTD plus some maximum tolerable delay time.

Fig. 7 shows the unified QoS factor of all frame types in this experimental system. The function has two parameters: early arrival defines the time the frame data arrives before TTD; timeout defines when the frame data arrives too late to be decoded. These two values are chosen according to application situations. A larger early arrival time may result in larger client-side buffer requirements and a larger timeout may cause delays during playback. On the other hand, smaller values will also affect delivered video quality.

In the described experimental system, some content information is also taken into consideration. From the viewpoint of sampling theory, we can see that frames in fast motion sequences should be preserved with higher precedence than those in slow motion sequences during the frame dropping process. It is also known that those frames containing more motion information will normally use more bits than those frames containing less motion information when they are predictive-coded. As a result, the QoS factor of node n is defined as $V(n, t) = U(t) * \text{coded_size}(n)$. However, its effects are very limited when video bitstreams are CBR coded.

System Implementation

In this described example, the prototype was implemented as a WWW service extension to MS Internet Information Server running on MS Windows NT. The adaptation application runs as an ISAPI extension on IIS. Video data is

processed and streamed in real-time from an original source through a standard HTTP protocol stack provided by IIS. A bandwidth-limitation software pipe was used as a simple emulation of network bandwidth. Adapted video data are firstly sent through this pipe before IIS sends it out. Parameters such as emulation bandwidth and optimizer search steps are all sent to the server as request parameters. Several popular client applications that support playback of MPEG I video have been successfully tested using HTTP streaming including Windows Media Player and QuickTime Player.

Experimental Results

We tested the implementation using both CBR and VBR bitstreams. Figs. 8 and 9 show some of these results. The 3Mbps VBR MPEG I bitstream is 320x240x30 fps and is two-pass-coded with minimal bitrate at 1Mbps and maximal bitrate at 4Mbps. Its GOP structure is 1I5P3B. The 1.2Mbps CBR bitstream is 320x240x29.97 fps and is one pass coded. Its GOP structure is 1I3P3B.

Fig. 8 shows delivered frame sizes of the first 200 frames of these two clips at different bandwidths. The curves show the average frames size over a window of 60 frames. One can clearly see which frames are escape-coded during the adaptive delivery process. The delivery bitrate of the 3Mbps VBR bitstream is also smoothed because we assumed fixed delivery bandwidth. Fig. 9 shows how the statistics of frame types change when the delivery bandwidth changes.

From these results one can see that bandwidths of the delivered streams are successfully reduced and smoothed. Thus, playback of these video bitstreams is possible even when the network bandwidth is far narrower than that the source

1 bitstreams demand and when the client side buffer size is limited. From Fig. 9 one
2 can see that frame data are preserved in the order of importance as preferred. This
3 simple adaptation scheme can have its limitation too. For example, the frame rates
4 of adaptation results are still not ideally adjusted because the structures of the
5 MPEG I video sequence are fixed and the QoS factor definition is not optimal.

7 Exemplary Application Programming Interfaces

8 Appearing below are a collection of exemplary application programming
9 interfaces (APIs) that can be utilized to implement embodiments of the system
10 described above.

```

12     template <class cYourEdge> class TcEdge
13     {
14     public:
15         TcEdge(int argnFromNodeID, int argnToNodeID):           // The
16         constructor      nFromNodeID(argnFromNodeID), nToNodeID(argnToNodeID),
17                         pNextOutEdge(0), pNextInEdge(0)
18                         {};
19         virtual ~TcEdge(){};
20
21         int nFromNodeID;           // Which node it is from
22         int nToNodeID;             // Which node it goes to
23         cYourEdge *pNextOutEdge;   // This is the outgoing edge link list
24         cYourEdge *pNextInEdge;   // This is the incoming edge link list
25     };
26
27     template <class TcEdge> class TcNode
28     {
29     public:
30         TcNode();
31         virtual ~TcNode();
32
33         TcEdge *pOutEdge;          // This is the outgoing edge link list
34         TcEdge *pInEdge;          // This is the incoming edge link list
35     };
36
37     template <class TcNode, class TcEdge> class TcGraph
38     {
39     public:
40         TcGraph();
41         virtual ~TcGraph();
42
43         virtual bool CleanUp();    // clean up all data in graph

```

```

1      virtual int AddNode(int nNodeIDPrefered, bool bPreferSmall=true);
2      // Add node to graph. You may specify a preferred nodeID or -1 for
3      // automatic assignments. Allocation will from small to larger if
4      // bPreferSmall is true
5      // return new Node ID on success or the preferred ID if exist, -1
6      on fail
7      virtual int AddNode(TcNode *pNode, int nNodeIDPrefered,
8      bool bPreferSmall=true);
9      // Add node and assign the node pointer
10     // return new Node ID on success or the preferred ID if exist, -1 on fail
11
12     virtual int DeleteNode(int nNodeID);
13     // return released Node ID on success, -1 on fail
14
15     virtual unsigned int GetNumNodes();
16     // return the number of nodes on the graph
17
18     virtual TcNode *GetNodePointer(int nNodeID);
19     virtual TcNode *SetNodePointer(TcNode *pNode, int nNodeID);
20     // Get and Set the node pointer
21
22     virtual TcEdge *AddEdge(int nFromNodeID, int nToNodeID);
23     // Add an edge from id1 to id2
24     // Return the edge pointer on success or NULL on failure
25     virtual int DeleteEdge(int nFromNodeID, int nToNodeID);
26     // Delete the edge from ID1 to ID2
27     // return the fromid on success and -1 on failure
28
29     virtual TcEdge *GetEdgePointer(int nFromNodeID, int nToNodeID);
30     // Get the edge pointer (from , to)
31
32     virtual bool StartNodeEnumeration();
33     // Start enumeration of existing nodes
34     // return true if it's ok
35     virtual int EnumerateNode();
36     // Get current enumerated node id and advance to next id
37     // return current node id on success, return -1 on end of
38     enumeration
39
40     virtual unsigned int PreRequireNodeBuffer(unsigned int
41     requiredsize);
42     // Used to pre-allocate node buffer, for better buffer management
43     only
44
45     private:
46     virtual bool AdjustNodeBuffer(int nMaxNodeID);
47     // adjust the node pointer buffer to hold new IDs
48     // return true if succeeded
49     // return false if failed or input parameter is not ok
50     // throw cGraphError with fatal errors
51
52     virtual int AllocNodeID(int nPreferedNodeID, bool
53     bPreferSmall=true);
54     virtual int FreeNodeID(int nNodeID);
55
56     // the following two is used internally for node buffer management
57     virtual void MoveMinFreeIDPointer();
58     virtual void MoveMaxFreeIDPointer();
59
60     };

```



```

1  //////////////////////////////////////
2  class cACRESEdge : public TcEdge<cACRESEdge>
3  {
4  public:
5      cACRESEdge(int argnFromNodeID, int argnToNodeID); // constructor
6      virtual ~cACRESEdge();
7
8      virtual bool UpdateStatus(ACRES_NODE_STATUS statusInNode);
9      // update the edge's status, check if it is ignition edge
10     // return true if it is.
11
12     // dump my self to a file
13     virtual bool WriteFile(FILE *fp);
14     // reinit my self from a file
15     virtual bool ReadFile(FILE *fp);
16
17     ACRES_EDGE_OBJ_TYPES objType;
18     // object type, needed by some application to distinguish
19     // between different classes (RTTI)
20
21     ACRES_EDGE_TYPE eType;
22     ACRES_EDGE_STATUS eStatus; // this is not used for pure route edges
23     ACRES_ROUTE_STATUS route; // this is not used for dependency edges
24     int eCondition; // should be const after been created;
25 };
26
27 class cACRESNode : public TcNode<cACRESEdge>
28 {
29 public:
30     cACRESNode();
31     cACRESNode(int hostID);
32     virtual ~cACRESNode();
33
34     virtual int CheckCondition();
35     // Check if the node can be put into activable one
36     // return
37     // > 0 active
38     // = 0 inactive
39     // < 0 permanent inactive because of cold_forever edge and
40     condition, // thus should be put into SKIPPED ASAP
41
42     virtual int UpdateOutEdgeStatus(bool includeRoute);
43     virtual int ResetOutEdgeStatus(bool includeRoute);
44     // return the number of changes made
45
46     // dump my self to a file
47     virtual bool WriteFile(FILE *fp);
48     // reinit my self from a file
49     virtual bool ReadFile(FILE *fp);
50
51     ACRES_NODE_OBJ_TYPES objType;
52     // node type, needed by some application to distinguish between
53     // different classes (RTTI)
54
55     ACRES_NODE_STATUS nStatus;
56     ACRES_NODE_CONDITION nCondition;
57     ACRES_NODE_BEHAVIOR behavior;
58
59     int nLayer; // required by ACRES validator

```



```

1      virtual int AllocHostID(int nPreferredNodeID, bool
bPreferSmall=true);
      virtual int FreeHostID(int nNodeID);
2
      virtual void MoveMinFreeHostIDPointer();
      virtual void MoveMaxFreeHostIDPointer();
3
    };
4
    class cACRESHost
5    {
    friend cACRESHost * cACRES::SetHost(cACRESHost *pHost, int hostID);
6    // this function will assign the pACRES and nHostID for this host object
    public:
7      cACRESHost() : pACRES(NULL), nHostID(-1),
objType(_ACRES_HOSTOBJ_BASIC)
8      {};
      ~cACRESHost(){};
9
      virtual bool Reset(void) {return false;};
      virtual cACRESNode * CreateNode()=0;
10     // The user application call this to create the content node
      virtual bool GetNodeContent(int nNodeID, ACRESContent *pCont)
11     {return false;};
      // Get the content represented by the Node
      virtual bool NotifyDeliveredNodeContent(int nNodeID) {return
12     false;};
      // The Decision engine call this if the node is delivered.
13
      virtual double GetNodeValue(int nNodeID, double nTime) {return 0;};
      virtual double GetNodeCost(int nNodeID, ACRES_NODE_STATUS
14     tostatus,
      double nTime) {return 0;};
      virtual double GetNodeResource(int nNodeID, ACRES_RESOURCE_TYPE
15     nResource)
      {return 0;};
      virtual double CheckPreferece(int nNodeID, void *pPref) {return
16     0;};
      // check compatiblilty, preference, profile, ...
17
      virtual bool GetTimePreference(int nNodeID, double *pTimePref)
18     {return false;};
19
      // time base of hosted contents, may only be useful for dynamic
      contents
      virtual double SetTimeBase(double timebase){return 0;};
      virtual double GetTimeBase() {return 0;};
20
      // dump my self to a file
      virtual bool WriteFile(FILE *fp) {return true;};
      // reinit my self from a file
      virtual bool ReadFile(FILE *fp) {return true;};
21
      // Add more common Host API here
22
      //////////// OBJ TYPE ////////////
23
      ACRES_HOSTOBJ_TYPES objType; // my RTTI
24
    protected:
25     cACRES *pACRES; // On which ACRES I am

```

```

1      int nHostID;          // My ID
2
3      ///////////////////////////////////////////////////
4      class cDecision
5      {
6      public:
7          cDecision(cTimeTicker *pTicker);
8          // the decision engine gets it's system clocks
9          ~cDecision();
10
11          virtual bool AssignContent(cACRES *pACRES);
12          // Give it our content
13
14          virtual bool AssignOutput(cOUTPUT *pOUT);
15          // Setup the output API
16
17          virtual bool setResource(ACRES_RESOURCE_TYPE restype,
18          cResourceModel *pRes);
19          // Setup the resource model
20
21          virtual bool SetLookForward(unsigned int nlookres,
22          unsigned int nlookdecision);
23          // Set the search parameter
24
25          virtual unsigned int GetSystemClock() {return sysclock-
>GetTick();};
26
27          virtual bool StartUp();
28          virtual int RunEngine(int *delivered);
29          // call it every clock tick
30          // return the number of delivered nodes in int *delivered
31          // and the number of status changes during the run
32
33          virtual bool ContentPending();
34          // check if there are still deliverable contents left
35
36      protected:
37          virtual bool Cleanup(bool all);
38
39          virtual int Qualificator();
40          // Update node statuses and prepare the candidate set
41
42          virtual int Planner();
43          // Run the decision procedure
44
45          virtual int Deliverer();
46          // Call the output API
47
48          virtual double CostEffective(double value, double cost,
49          double timenow, double timepref);
50          // calculating the value of content
51
52          // Implementation of optimized searching algorithm
53          //
54          virtual bool InitFuturePlanner(unsigned int maxPredicts);
55          virtual double FuturePlannerCheckDelay(double res);
56          virtual bool FuturePlannerAllocateResource(int nID, cACRESHost
57          *pHost,

```

```

double &delay);
virtual bool FuturePlannerFreeResource(int nID, cACRESHost *pHost,
double &delay);
virtual int FuturePlannerExpand(int level, int maxlevel,
int nTryNodeID, cACRESNode *pTryNode,
cDENodeList *pFire, cACRESEdge *pEdge,
double &totalbenifit, double &totalvalue, double
&totalcost);
virtual int FuturePlanner(int level, int maxlevel,
double &totalbenifit, double &totalvalue, double &totalcost,
int &idFrom);
virtual bool DeInitFuturePlanner();

// used for resource management during the decision process
virtual double ResourceManager(bool init);
virtual double CheckResource(ACRES_RESOURCE_TYPE restype);
virtual double LookForwardResource(ACRES_RESOURCE_TYPE restype,
unsigned int systick);
virtual bool ReserveResource(int nID, cACRESNode *pNode, double
*delay);

// These two routines are used to maintain the candidate sets
which
// are implemented as link lists
virtual cDENodeList* AddTo(cDENodeList* &rpList, int nid,
cACRESNode *pNode, int fromid);
virtual cDENodeList* DeleteFrom(cDENodeList* &rpList,
cDENodeList *pElem, bool drop=true);

virtual int UpdateRelatedNodeStatus(int nid, cACRESNode *pNode,
bool updateEdge, int fromid);
// update the status of related nodes if my status is changed

virtual int UpdateRouteToWarmLine(int nid, cACRESNode *pNode);
// Update WarmLine nodes from a FireLine Node
// return the possible usable routes of the FireLineNode;
// return -1 for errors

virtual int CheckRouteToWarmLine(int nid, cACRESNode *pNode);
// check if this node can goto warmline or potentially goto
warmline
// return -1 for errors
// return 0 for not be able to reach any inactive or activable
nodes
// return 1 for yes

virtual cDENodeList * FindInList(int nid, cDENodeList *pList);
// find the node in a list sorted increasely by Layer and NodeID

int SkipValuelessRouteNodes(int nid, cACRESNode *pNode, int
&numskip);
// return the number of valuable/branchable route nodes
};

```

099499 " 6545650

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Conclusion

Methods and systems that provide a framework for generic adaptive multimedia content delivery have been described. The framework features an abstract content model and an abstract adaptation decision engine that can make adaptive delivery plans without knowing much of the physical details of actual content. The capabilities of the framework have been demonstrated with an application of adaptive video streaming. Experimental results further show that the proposed framework is effective and efficient in adaptive delivery of contents under variable network conditions. The described architecture can be easily extended to have much stronger capabilities.

Although the invention has been described in language specific to structural features and/or methodological steps, it is to be understood that the invention defined in the appended claims is not necessarily limited to the specific features or steps described. Rather, the specific features and steps are disclosed as preferred forms of implementing the claimed invention.